

---

Citation:

Altahhan, A (2018) TD(0)-Replay: An Efficient Model-Free Planning with full Replay. Proceedings of International Joint Conference on Neural Networks. ISSN 2161-4393 DOI: <https://doi.org/10.1109/IJCNN.2018.8489300>

Link to Leeds Beckett Repository record:

<https://eprints.leedsbeckett.ac.uk/id/eprint/4849/>

Document Version:

Article (Accepted Version)

---

The aim of the Leeds Beckett Repository is to provide open access to our research, as required by funder policies and permitted by publishers and copyright law.

The Leeds Beckett repository holds a wide range of publications, each of which has been checked for copyright and the relevant embargo period has been applied by the Research Services team.

We operate on a standard take-down policy. If you are the author or publisher of an output and you would like it removed from the repository, please [contact us](#) and we will investigate on a case-by-case basis.

Each thesis in the repository has been cleared where necessary by the author for third party copyright. If you would like a thesis to be removed from the repository or believe there is an issue with copyright, please contact us on [openaccess@leedsbeckett.ac.uk](mailto:openaccess@leedsbeckett.ac.uk) and we will investigate on a case-by-case basis.

# *TD(0)-Replay: An Efficient Model-Free Planning with full Replay*

Abdulrahman Altahhan  
School of Computing and Creative Technology  
Leeds Beckett University  
Leeds, UK  
a.altahhan@leedsbeckett.ac.uk

**Abstract**— In this paper we present a two novel reinforcement learning methods that allow for full replay of all past experience in every step of a reinforcement learning agent life with minimal overhead. In particular, we show how to deduce an equivalent efficient backward view of replaying the full past experience online using TD(0) error. We emphasise the already established link between replaying and planning in our algorithm design by comparing it with an Extensive Linear Dyna Planning algorithm, and we show that our method can outperform this expensive form of planning methods. We test the new method, which we called TD(0)-Replay, on two different domains problems; Dyna Maze to test its planning capabilities, and Random Walk to test its prediction capabilities. We compare TD(0)-Replay with TD( $\lambda$ ) for benchmarking and we show that our method outperform this traditional RL method as well. We also show that our method when combined with weight reinitialisation turns into especially effective form of planning.

**Keywords**—component, formatting, style, styling, insert (key words)

## I. INTRODUCTION

In Reinforcement Learning, replaying past experience has been shown to have an important and definite role in reaching an optimal or close-to-optimal policy. Replay becomes even more important when dealing with experience that is difficult or expensive to simulate or when the trajectory of available experience is very limited. Especially in applications that requires real world interaction, it is difficult, and simply undesirable, to have to repeated the experience physically. Instead replaying the experience in the mind of the agent and learning from it becomes a natural and important method of learning. In essence, such complex task that needs rich imagination is going to be inevitably computationally expensive. Yet, in online learning, allowing the agent to *fully* repeat *all* past experience that has been done *so far* is even more demanding. However, this full repetition if achieved efficiently, can provide the agent with a very powerful learning mechanism that boosts its performance and allows it to maximise the so far lived experience in a way that has not been done before.

In this paper we will provide a new method that allows the agent to achieve exactly the above. I.e. to allow the agent to fully replay all past experience (in its head rather than physically), which will allow it, in turn, to quickly optimise its value function prediction as well as improve its policy. We will show that this

can be achieved with a reasonable computational expense that makes the proposed method, called TD-Replay, a very attractive method for the above mentioned situations which requires maximisation of so far experience without having to physically repeat it. [1] and [2] for example studied the effect of replaying and they have shown that the agent can boost its experience when using replaying. The original interpretation for experience replay is that the samples will be presented for the agent as a new set of samples [1]. While in [3] the agent suffices by re-updating its weight as well as its value function estimation either fully or partially. At the same time, [4] studied the effect of replaying from planning perspective and they showed that replaying can be looked at as planning by looking into the past (instead of future) and they showed that their TD(0) replay algorithm is equivalent to the Planning with the linear Dyna model Algorithm. We will follow a similar approach as in [4] and [5], however we will develop a different algorithm that has its own update rules and mechanism that is different from the presented algorithms but has some resemblance in terms of the form of updates.

In our work we will introduce an efficient algorithm with forward view that depends on the mechanism provided by [4] and [5]. True Online TD [5] allows the learning process to be repeated for all  $k = 1 \dots t$  but it does not utilise replay, it assumes that the agent will always reinitialise its weights to the same initial values in every time step  $t$ . On the other hand, [4] utilises replaying on the level of targets only. Replaying TD(0) updates Algorithm for example (and its more efficient equivalent Planning with Dyna Algorithm) assumes that the agent starts from the same initial weights at every set of imaginary experience, only its  $U_i$  targets are changing according to the latest weights coming from the real time step  $t$ . Our algorithm, which we call TD(0)-Replay, assumes that both the targets as well as the initial weights are changing in every set of imaginary updates  $k = 1 \dots t$ , making it more vigilant and adaptive to changes in the environment.

TD(0)-Replay will depends on the usual TD error and will not reinitialise the weights in each time step  $t$ , instead it will assume that the starting weights of time step  $t$  are those obtained after updating the weights in time step  $t - 1$ , but will assume that the agent is going to replay all of its past experience for each imaginary time step  $k$ , where  $k = 1 \dots t$ .

## II. TD-REPLAY FORWARDVIEW

### A. TD-Replay update rules at time step $t$

In this section we will develop the updates rules for our TD-Replay method on the basis of a one-layer neural network model (linear model). At time step  $t$  TD(0) update for a linear model is given as

$$\delta_t = R_{t+1} + \gamma \theta_t^\top \phi_{t+1} - \theta_t^\top \phi_t \quad (1)$$

$$\theta_{t+1} = \theta_t + \alpha_t \delta_t \phi_t^\top \quad (2)$$

Where:  $\delta_t$  is the temporal Difference error,  $R_{t+1}$  is the reward signal,  $\gamma$  is a discount factor,  $\phi_t^\top$  is the transpose of feature vector  $\phi_t$  obtained through current state  $S_t$ ,  $\theta_t^\top$  is the transpose of the weight vector  $\theta_t$  and  $\alpha_t$  is a learning step; all varies according to time step  $t$ .

In order to replay previous experience and assuming that the agent is at time step  $t$ , in this case the TD(0) error for a time step  $k$ :  $0 \leq k < t$  is given by

$$\theta_{k+1} = \theta_k + \alpha_k (R_{k+1} + \gamma \theta_k^\top \phi_{k+1} - \theta_k^\top \phi_k) \phi_k \quad (3)$$

Or

$$\theta_{k+1} = [I + \alpha_k \phi_k (\gamma \phi_{k+1}^\top - \phi_k^\top)] \theta_k + \alpha_k R_{k+1} \phi_k \quad (4)$$

Hence by renaming we have:

$$\theta_{k+1} = A_k \theta_k + B_k \quad (5)$$

$$A_k := [I_{n \times n} + \alpha_k \phi_k (\gamma \phi_{k+1}^\top - \phi_k^\top)] \quad (6)$$

$$B_k := \alpha_k \phi_k R_{k+1} \quad (7)$$

Where  $n$  is the features dimension i.e.  $|\phi_k| = n$ . Of course since we are using a linear model we have also  $|\theta_k| = n$ .  $A_k$  is an  $n \times n$  squared matrix, while  $B_k$  is  $n \times 1$  vector,  $\top$  is the transpose symbol.

### B. Cumulative update rules at time steps $t = k + 1$

If we allow the agent to replay all so far experience in every step, we get the following formulae:

$$\begin{aligned} t = 1 \quad k = 0 \quad & \theta_1^1 = A_0 \theta_{init} + B_0 \\ t = 2 \quad k = 0 \quad & \theta_1^2 = A_0 \theta_1^1 + B_0 \\ & k = 1 \quad \theta_2^2 = A_1 \theta_1^2 + B_1 \\ t = 3 \quad k = 0 \quad & \theta_1^3 = A_0 \theta_2^2 + B_0 \\ & k = 1 \quad \theta_2^3 = A_1 \theta_1^3 + B_1 \\ & k = 2 \quad \theta_3^3 = A_2 \theta_2^3 + B_2 \\ & \dots \end{aligned} \quad (8)$$

By convention, since the last imaginary step is  $k = t - 1$  for  $t$ , (and the last update is going to be  $\theta_{k+1}^t = \theta_t^t$ ), we define  $\theta_t := \theta_t^t$ . The above will give the following set of formula for  $\theta_t^t$ .

$$\begin{aligned} \theta_1 &= A_0 \theta_0^0 + B_0 \\ \theta_2 &= A_1 A_0 \theta_1 + A_1 B_0 + B_1 \\ \theta_3 &= A_2 A_1 A_0 \theta_2 + A_2 A_1 B_0 + A_2 B_1 + B_2 \end{aligned}$$

...

$$\theta_{t+1} = A_t \dots A_0 \theta_t + A_t \dots A_1 B_0 + \dots + A_t B_{t-1} + B_t \quad (9)$$

By defining

$$A_t^i := A_t \dots A_i \quad (10)$$

$$A_t^{t+1} := I_{n \times n} \quad (11)$$

The previous equations can be written through induction as

$$\theta_{t+1} = A_t^0 \theta_t + A_t^1 B_0 + \dots + A_t^t B_{t-1} + A_t^{t+1} B_t \quad (12)$$

Where:  $A_t := [I_{n \times n} + \alpha_t \phi_t (\gamma \phi_{t+1}^\top - \phi_t^\top)]$

$$\theta_{t+1} = A_t^0 \theta_t + \sum_{i=0}^t A_t^{i+1} B_i \quad (13)$$

Hence by defining the eligibility trace  $e_t$  and eligibility matrix  $\hat{e}_t$  as:

$$\hat{e}_t := A_t^0 \quad (14)$$

$$e_t := \sum_{i=0}^t A_t^{i+1} B_i \quad (15)$$

The updates rules for the learning weights can be written as:

$$\theta_{t+1} = \hat{e}_t \theta_t + e_t \quad (16)$$

### C. Incremental cumulative update rules at time steps $t = k + 1$

Let us now deduce incremental rules for the eligibility trace  $e_t$  and eligibility matrix  $\hat{e}_t$ .

As for the eligibility matrix  $\hat{e}_t$ , by induction we have:

$$\hat{e}_t = A_t^0 = A_t A_{t-1} \dots A_0 = A_t A_{t-1}^0 = A_t \hat{e}_{t-1}.$$

While for the eligibility trace  $e_t$ , we have

$$e_t = \sum_{i=0}^t A_t^{i+1} B_i = A_t \sum_{i=0}^{t-1} A_{t-1}^{i+1} B_i + B_t = A_t e_{t-1} + B_t, \text{ since we have } e_{t-1} = \sum_{i=0}^{t-1} A_{t-1}^{i+1} B_i.$$

Hence our TD(0)-Replay algorithm can be written in the following order:

$$A_t = [I_{n \times n} + \alpha_t \phi_t (\gamma \phi_{t+1}^\top - \phi_t^\top)] \quad (17)$$

$$B_t = \alpha_t \phi_t R_{t+1} \quad (18)$$

$$e_t = A_t e_{t-1} + B_t \quad (19)$$

$$\hat{e}_t = A_t \hat{e}_{t-1} \quad (20)$$

$$\theta_{t+1} = \hat{e}_t \theta_t + e_t \quad (21)$$

Where we have  $A_0 = \hat{e}_0 = I_{n \times n}$ ,  $B_0 = \alpha_0 \phi_0 R_1$ ,  $e_0 = 0_{n \times 1}$ . It should be noted that the algorithm uses just current time step information to apply a full replay of all past experience hence its significance lies in this particular characteristic.

#### D. Efficient Form of TD-Replay Forward

TD-Replay in its previous form can be made more efficient by unpacking  $A_t$  in the updates and replacing matrix multiplications in (19) and (20) with matrix to vector multiplication.

As for the calculations of  $e_t$  we have  $e_t = A_t e_{t-1} + B_t$ , hence  $e_t = e_{t-1} + \alpha_t \phi_t [(\gamma \phi_{t+1} - \phi_t)^\top e_{t-1}] + B_t$ . Therefore the calculations of involves calculating the term  $(\gamma \phi_{t+1} - \phi_t)^\top e_{t-1}$  is a scalar (multiplying a vector  $e_{t-1}$  by a vector transpose) which is more efficient than multiplying a squared matrix  $A_t$  and a vector  $e_{t-1}$ .

As for  $\hat{e}_t$  we have  $\hat{e}_t = A_t \hat{e}_{t-1}$ , hence  $\hat{e}_t = \hat{e}_{t-1} + \alpha_t \phi_t [(\gamma \phi_{t+1} - \phi_t)^\top \hat{e}_{t-1}]$ . It should be noted that  $[(\gamma \phi_{t+1} - \phi_t)^\top \hat{e}_{t-1}]$  is a multiplication of a squared matrix  $\hat{e}_{t-1}$  by a vector  $(\gamma \phi_{t+1} - \phi_t)^\top$  and is more efficient than multiplying two squared matrices  $A_t$  and  $\hat{e}_{t-1}$  as before. The complexity still lies within this calculation which is of  $O(n^2)$  for space and time. The final algorithm can be written as:

$$e_t = e_{t-1} + \alpha_t \phi_t [(\gamma \phi_{t+1} - \phi_t)^\top e_{t-1} + R_{t+1}] \quad (22)$$

$$\hat{e}_t = \hat{e}_{t-1} + \alpha_t \phi_t [(\gamma \phi_{t+1} - \phi_t)^\top \hat{e}_{t-1}] \quad (23)$$

$$\theta_{t+1} = \hat{e}_t \theta_t + e_t \quad (24)$$

Where we have  $e_0 = 0_{n \times 1}$ ,  $\hat{e}_0 = I_{n \times n}$

#### E. TD(0)-Replay Forward Algorithm

Formula (22) -(24) define a set of update rules for an agent to be able to predict the value function for a specific task in some environment. The algorithm is given below.

---

##### Algorithm 1 TD(0)-Replay: Value Function Prediction

---

INPUT:  $\alpha, \gamma, \theta_{init}$

$\theta \leftarrow \theta_{init}$

Loop (over episodes):

Obtain initial  $S, \phi$

$\hat{e} \leftarrow I_{n \times n}$ ,  $e \leftarrow 0_{n \times 1}$

While (terminal state has not been reached), do:

act according to the policy

observe next reward  $R = R_{t+1}$ , next state  $\hat{S} = S_{t+1}$

and its features  $\hat{\phi} = \phi_{t+1}$

$\alpha \leftarrow \ell(\alpha)$

$\theta \leftarrow \theta + \alpha \phi [(\gamma \hat{\phi} - \phi)^\top \theta + R]$  ( $\delta$  update)

$e \leftarrow e + \alpha \phi [(\gamma \hat{\phi} - \phi)^\top e + R]$  (Re-Playing/Re-Planning)

$\hat{e} \leftarrow \hat{e} + \alpha \phi [(\gamma \hat{\phi} - \phi)^\top \hat{e}]$

$\theta \leftarrow \hat{e} \theta + e$

$\phi \leftarrow \hat{\phi}$

---

It should be noted that the agent is applying its current update through the normal delta update then it does a full replay again from first step up until the last step. Therefore, effectively the agent takes into consideration current step information (rewards and feature), updates the weights, replay all past experience and updates the weights accordingly, then finally re-update the weights of current step again according to the replay. This is

more effective than replaying up until the step before the last. In the above algorithm the learning rate  $\alpha$  is assigned at every time steps  $t$  according to  $\ell(\alpha)$  which can be any scheme that reduces  $\alpha$  (annealing for example). In practices if  $\alpha$  is chosen to be small enough then it can be left without updating it [3].

#### F. Comparing TD(0)-Replay with Other Replay Algorithms

Clearly from a formative perspective the updates rules are different than those presented in From [4]. From a fundamental inner working mechanism, the main differences between ‘TD(0)-Replay’ and ‘Replaying TD(0)’ Algorithms [4] (and to some extent even the True online TD algorithms [5]) are in three essential characteristics. The first, is that the dynamics are different since in [4] the agent reflects back and uses the latest weights that resulted from previous step  $t$  in calculating the targets  $U_k$  in all consequent steps  $k = 0 \dots t$ . While, TD(0)-Replay uses the weights  $\theta_t^t$  from past experience  $t$  as an initial weights in step  $k = 0$ , then it lets the consequent updates specify targets  $U_k$  that has been calculated using the latest  $\theta_k^{t+1}$  in consequent steps  $k = 1 \dots t$ . Secondly, the term  $A_t$  is different since we have  $A_t = [I_{n \times n} + \alpha_t \phi_t (\gamma \phi_{t+1} - \phi_t)^\top]$  instead of  $A_t = [I_{n \times n} - \alpha_t \phi_t (\phi_t)^\top]$ . Lastly, the broader concept of weight updates with no reinitialisation makes our approach more general. In other words for this form of TD-Replay we do not reinitialise the weights in the start of each imaginary set of updates unlike [4][5] and [10]. Later we will relax this assumption to obtain a new algorithm.

On the other hands, both TD(0)-Replay and Replaying TD(0), replay fully all previous experience in every steps  $t$ . Also both do not require storing any previous states or weights, hence are efficient. TD-Replay only requires storing the eligibility squared matrix  $\hat{e}$  and  $e$  (equivalent requirements for Replaying TD(0) updates is to store square matrix  $F$  and  $b$ ). The complexities of both our algorithm and their algorithms for storage and computation are of  $O(n^2)$  in the worst case regardless of the number of steps.

#### G. TD(0)-Replay and Re-Planning

From another perspective, since we are updating the weights then reflecting back on the agent past experience to learn from it, it is only reasonable to consider our algorithm as a re-planning algorithm as well.  $\hat{e}$  can be considered to be trying to establish a prediction model for the feature vector that is capable of being changed according to the difference between the two consequent vectors  $(\gamma \hat{\phi} - \phi)^\top$  instead of predicting the next feature vector as in [3]. This makes sense since the TD(0)-Replay algorithm does not need to know the next feature, instead it needs to know how the features are changing. Similarly, it needs to know how the reward function will change through  $e$  in order to come up with a planning scheme for the future replays of past experience according to the latest up to date weights and targets. It should be noted that TD(0)-Replay is different than an algorithm that utilises the residual gradient since the term  $(\gamma \hat{\phi} - \phi)^\top \hat{e}$  is multiplied by  $\alpha \phi$  not with  $(\gamma \hat{\phi} - \phi)^\top$ .

#### H. TD(0)-Replay with Reinitialisation

If we changed the mechanism of TD(0)-Replay so that it will reinitialise its weights in each imaginary step  $k = 0$ , to  $\theta_0^t = \theta_{init} \forall t$  then the update rules (17)-(20) stay the same, update rule (21) becomes

$$\theta_{t+1} = \hat{e}_t \theta_{init} + e_t \quad (25)$$

And the next update is going to be

$$\theta_{t+2} = \hat{e}_{t+1} \theta_{init} + e_{t+1}$$

In this case, since  $\theta_{init}$  is fixed then we can make the updates rules much more efficient since we can store  $\hat{e}_t \theta_{init}$  instead of  $\hat{e}_t$ . Hence, the update rule (20) can be changed into:

$$\hat{e}_t \theta_{init} = A_t \hat{e}_{t-1} \theta_{init} \quad (26)$$

By defining  $\hat{\theta}_t$  as a vector

$$\hat{\theta}_t = \hat{e}_t \theta_{init} \quad (27)$$

The update (28) can be written as:

$$\hat{\theta}_t = A_t \hat{\theta}_{t-1} \quad (28)$$

Hence, these updates along with (22) define a new algorithm that can be written efficiently as

$$e_t = e_{t-1} + \alpha_t \phi_t [(\gamma \phi_{t+1} - \phi_t)^\top e_{t-1} + R_{t+1}] \quad (29)$$

$$\hat{\theta}_t = \hat{\theta}_{t-1} + \alpha_t \phi_t [(\gamma \phi_{t+1} - \phi_t)^\top \hat{\theta}_{t-1}] \quad (30)$$

$$\theta_{t+1} = \hat{\theta}_t + e_t \quad (31)$$

Where we have  $e_0 = 0_{n \times 1}$ ,  $\hat{\theta}_0 = \theta_{init}$  all as vectors. The algorithm is shown below.

---

**Algorithm 2** *TD(0)-Replay( $\theta_0$ ) with Reinitialisation: Value Function Prediction(Policy Evaluation)*

---

INPUT:  $\alpha, \gamma, \theta_{init}$

$\theta \leftarrow \theta_{init}$

Loop (over episodes):

Obtain initial  $S, \phi$

$\hat{\theta} \leftarrow \theta, e \leftarrow 0_{n \times 1}$

While (terminal state has not been reached), do:

act according to the policy

observe next reward  $R = R_{t+1}$ , next state  $\hat{S} = S_{t+1}$

and its features  $\hat{\phi} = \phi_{t+1}$

$\alpha \leftarrow \ell(\alpha)$

$e \leftarrow e + \alpha \phi ((\gamma \hat{\phi} - \phi)^\top e + R)$  (replaying)

$\hat{\theta} \leftarrow \hat{\theta} + \alpha \phi ((\gamma \hat{\phi} - \phi)^\top \hat{\theta})$

$\theta \leftarrow \hat{\theta} + e$

$\phi \leftarrow \hat{\phi}$

---

One important issue to realise in this algorithm is the initialisation to  $\hat{\theta} \leftarrow \theta$  in each episode. This is important so that the updates of  $\theta$  are not lost. Also we have removed the ( $\delta$

update steps because it has no effect since the weights are effectively being reinitialised in each step. This new algorithm, regardless of the policy boosting update step, is doing a special type of efficient replay that resemble the one done in True Online TD but it uses one step backup rather than the full interim  $\lambda$ -returns. Moreover, it is not a special case or equivalent to True Online TD(0), since True Online TD(0) is equivalent to normal TD(0), while obviously this algorithm is doing more updates than TD(0). This algorithm can be considered a counter part of TD( $\lambda$ ) but it is analytical and is also truly online because it has an equivalent online algorithm.

Algorithms 1 and 2 have exact equivalent algorithms that are based on some form of repetitive TD updates. *TD(0)-Replay* can also be generalised by using the full interim  $\lambda$ -returns but this will be left for future work.

We will call Algorithm 2 *TD(0)-Replay( $\theta_0$ )* to recognise the fact that in each step it starts from the initial weight  $\theta_0 = \theta_{init}$ , while the original *TD(0)-Replay* will be called *TD(0)-Replay( $\theta_t$ )* to reflect the fact that in each time step it starts from the latest weights  $\theta_t$ .

#### I. Sarsa(0)-Replay Forward Algorithm

As for control, we can straightforwardly build control algorithm on the grounds of TF-Replay. The agent would need to learn a suitable policy; the policy can be deduced out of the predicted value function for the agent.

One way to build the policy learning model, based on the value prediction model, is as follows. A set of learning weights is provided for each action, given the set of actions are limited. The control features  $\psi_t$  are going to have a cardinality of  $|\psi_t| = |\mathcal{A}| \times n$  where  $n = |\phi_t|$  (the state feature size) and  $|\mathcal{A}|$  is the number of actions that an agent can take at any time step  $t$ . The agent also would have the same size for its weights when using a linear model i.e.  $|\mathcal{A}| \times n$  weights.

In each time step, the set of features  $\phi_t(a_t)$  in  $\psi_t$  (that corresponds to the current action  $a_t$ ) will be populated with the values of the state features, while the rest  $\phi_t(a_i): a_i \neq a_t$  in  $\psi_t$  will be simply populated with 0. Hence the learning takes place on the set of weights corresponding to the current actions since the rest of the features are going to be 0.

In order to deduce a suitable policy, the agent calculates the value function for each action and then picks the action with the highest value (most of the time, except for few times with small probability of  $\epsilon$  where the agent picks a random action). This type of policy, called  $\epsilon$ -greedy policy, is a common policy to be followed, other policies such as soft-max is also possible.

In this paper we will follow an  $\epsilon$ -greedy policy. By doing the above scheme a similar algorithm can be written for the agent in order to learn a suitable policy instead of learning only to predict the value function of its current policy. According to the policy improvement theorem this scheme of improving the policy by picking the max action value then updating the prediction accordingly will lead to convergence to an optimal policy in the case of a linear model that is being updated according to the TD error, with some extra conditions on the learning rate [6][7]. The policy improvement algorithm for

*Sarsa(0)-Replay* is given above. It should be noted that by convention  $\dot{\psi} \leftarrow 0$  if  $\dot{S}$  is terminal.

In order to improve the policy an agent can run either indefinitely or through a set of episodes that is specified in priori, where it stops when the learning slows down under a specific threshold.

For our comparisons we will choose a specific number of episodes and compare the Route Mean Squared Error (RMSE) or the total number of steps to for those episodes.

---

**Algorithm 3** *Sarsa(0)-Replay( $\theta_t$ ): Policy Improvement*

---

INPUT:  $\alpha, \gamma, \theta_{init}$   
 $\theta \leftarrow \theta_{init}$   
 Loop (over episodes):  
   Obtain initial  $S, \phi$   
   Select action  $A$  based on State  $S$   
    $\psi \leftarrow$  features corresponding to  $S, A$ ; ( $|\psi| = n \times |\mathcal{A}| = \mathcal{N}$ )  
    $\hat{e} \leftarrow I_{\mathcal{N} \times \mathcal{N}}, e \leftarrow \theta \leftarrow \mathbf{0}_{\mathcal{N} \times 1}$   
   While terminal state has not been reached, do:  
     take action  $A$ , observe next state  $\dot{S}$  and reward  $R$   
      $a \leftarrow \epsilon$ -greedy( $\arg\max \dot{Q}(A_i) \leftarrow \theta^T(a_i)\dot{\phi}$ ) ( $|\dot{Q}| = |\mathcal{A}|$ )  
      $\dot{\psi} \leftarrow$  features corresponding to  $\dot{S}, a$   
      $\alpha \leftarrow l(\alpha)$   
      $\theta \leftarrow \theta + \alpha \phi ((\gamma \dot{\psi} - \psi)^T \theta + R)$  ( $\delta$  update)  
      $e \leftarrow e + \alpha \psi ((\gamma \dot{\psi} - \psi)^T e + R)$  (re-planning)  
      $\hat{e} \leftarrow \hat{e} + \alpha \psi [(\gamma \dot{\psi} - \psi)^T \hat{e}]$   
      $\theta \leftarrow \hat{e} \theta + e$   
      $\psi \leftarrow \dot{\psi}; A \leftarrow \dot{A}$

---

### J. Dyna Full Planning Algorithm

In order to compare our algorithms objectively we will compare TD(0)-Replay with a special version of Dyna Planning which we call Dyna Full Re-Planning, where the agent regenerate (reimagine) fully all previous samples in every time step in order to better plan what to do with them in case it sees them in the future.

It seems reasonable that the concept of re-planning based on replaying past experience, works especially when the agent is expected to revisit some states due to its incompetence of reaching its terminal state or achieving its final goal (the case in lots of RL environment and tasks). The algorithm, shown above, is expensive because its complexity is going to be  $\sum_{t=1}^T t = O(T^2)$  rather than  $O(n^2)$ , where  $n$  is the feature size and  $T$  is the total number of visited states.  $T$  is normally  $\gg n$  especially at the start if learning. This extreme case of planning is conceived as the maximum performance any Dyna Re-Planning algorithm can achieve. We will compare this algorithms performance with ours to show the real planning capabilities of TD-Replay algorithms. Similar to the other algorithms, a policy improvement algorithm can be devised based on this above.

---

**Algorithm 4** *Dyna Full Re-Planning: Extreme Planning for Policy Evaluation (expensive; for comparison only)*

---

INPUT:  $\alpha, \gamma, \theta_{init}$   
 $\theta \leftarrow \theta_{init}$   
 Loop (over episodes):  
   Obtain initial  $S, \phi$   
    $F \leftarrow \mathbf{0}_{n \times n}, e \leftarrow \mathbf{0}_{n \times 1} t \leftarrow 1$   
   While (terminal state has not been reached), do:  
     act according to the policy observe next reward  $R, \dot{S}$  and  $\dot{\phi}$   
      $\theta = \theta + \alpha [R + \theta^T \dot{\phi} - \theta^T \phi] \phi$   
      $F \leftarrow F + \alpha [\gamma \dot{\phi} - F \phi] \phi^T$  (F is to predict next state  $\dot{\phi}$ )  
      $b \leftarrow b + \alpha (R - b^T \phi) \phi$  (b is to predict next reward  $R$ )  
      $\phi_t \leftarrow \phi$  (store visited state features, memory expensive)  
     For  $k \leftarrow 1$  to  $t$  (planning steps, computationally expensive)  
       based on revisiting past states  
       in order to obtain a better policy evaluation  
        $\dot{\phi} \leftarrow F \phi_k$   
        $R \leftarrow b^T \phi_k$   
        $\theta \leftarrow \theta + \alpha [R + \theta^T \dot{\phi}_k - \theta^T \phi_k] \phi_k$   
      $t \leftarrow t + 1$

---

### III. EXPERIMENTS DESIGN AND ALGORITHM TESTING

We have tested our algorithm on a two test beds. The first is Random Walk which is a Markov Reward Process (MRP) to test the TD(0)-Replay prediction algorithm. MRPs are useful tools to isolate the prediction problem from the policy improvement (control) problems. The idea behind it is to assign an action, in each step, based on a transition probability that represents the dynamics of the environment. The actions are generated due to this probability only, there is no decision making taking place and the policy is stochastic with fix probability. We will use a 6-state Random Walk environment, where the process starts from a middle state as in Fig. 1. The current state will be moved to the state in the left according to a probability of  $p$ , or to the right according to the probability  $1-p$ . Once the process reaches the final state to the right the process stops and the agent is rewarded  $+1$ , otherwise if the process reaches the final state to the left, the process stops and the agent is rewarded with 0. All other transitions have a 0 reward.

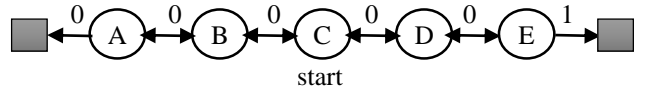


Fig. 1. Random Walk for 6 states. The true value for those states are their probability of reaching the far left terminal state assuming the agent start from each of them.

We have set  $\gamma = 1$  and we have used a very simple set of binary features that each represents a state. The features size is equal to the number of states. we have studies the effect of the learning rate for TD(0)-Replay in comparison with accumulate TD( $\lambda$ ) as a benchmark as well as with Dyna Full Re-Planning algorithms. Fig. 2 shows the results. It seems, due to its simplicity and effectiveness, the winner in this testbed is TD(0)-Replay( $\theta_0$ ), in this simple prediction domain.



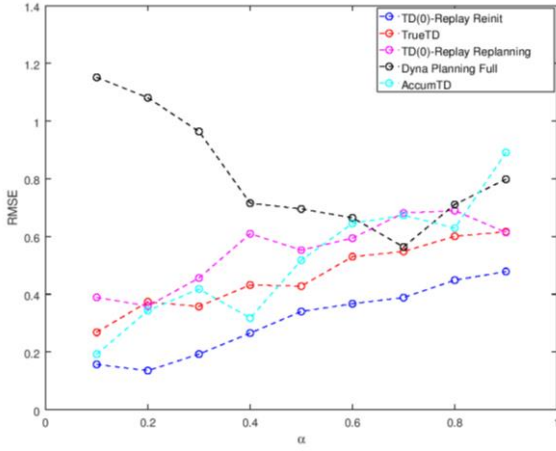


Fig. 2. Comparison of RMSE of different algorithms for the Random Walk problem and different learning rates. Clearly TD(0)-Replay( $\theta_0$ ) performed best due to its simplicity and effectiveness in this prediction problem, overcoming TD(0)-Replay( $\theta_0$ ) as well as True Online TD( $\lambda$ ,  $\lambda = 0.7$ ), algorithms.

The second test bed the traditional Dyna Maze environment with  $9 \times 6$  cells (states). The main goal of a Dyna Maze is for an agent to be able to learn a policy that allows it to reach a specific goal state where it starts from a specific initial state (cell) in each episode. The episode ends when the agent reaches the goal square. Plenty of obstacles have been placed in the way from the Start state to the goal state. The agent is rewarded with -1 in each step wasted before reaching the goal while the reward for reaching the goal state is 0. An example of a simulated agent represented as a red square is shown in Fig. 3.

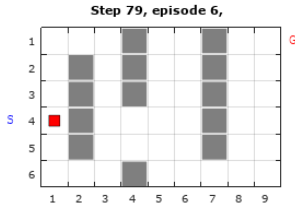


Fig. 3. Dyna Maze example of the small environment

Dyna Maze is traditionally used to test planning algorithms, we show here that *Sarsa(0)-Replay* exceeds the performance of other RL algorithms that involves planning. It actually competes head to head with very expensive re-planning algorithm, Dyna Full Re-Planning, that covers all previously visited states in each current state update. It should be noted that in the below experiments all results have been produced after averaging 20 runs. We have used binary features to represent the states by linearising the domain. We have 54 states each represented as a vector with 54 features, and we set  $\gamma = .95$ .

Fig. 4 shows the results for the Dyna Full Re-Planning algorithm mentioned before. It should be mentioned that this algorithm is included for comparison purpose only since it is expensive and can be impractical. Fig. 5 shows the comparison for *Sarsa(0)-Replay*( $\theta_t$ ).

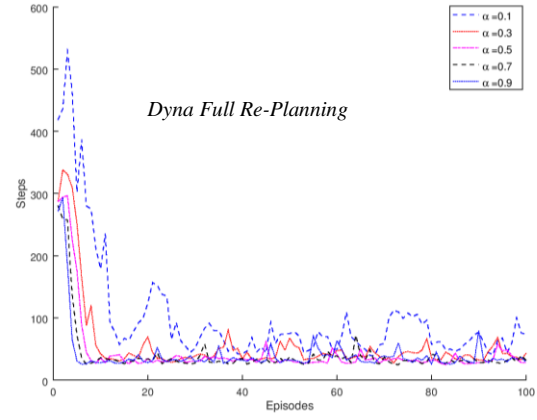


Fig. 4. Comparison of different learning steps for Dyna Full Re-Planning for control averaged over 20 runs for binary features. Clearly 0.9 is best and the algorithm converges very quickly after the 5<sup>th</sup> episode in all cases, lower alpha gave more stability for the algorithm. But the algorithm is very slow and inefficient (Complexity is  $O(T^2)$  in terms of time and memory,  $T$  is the total number of steps) because the agent resamples all of the so far visited states in every time step.

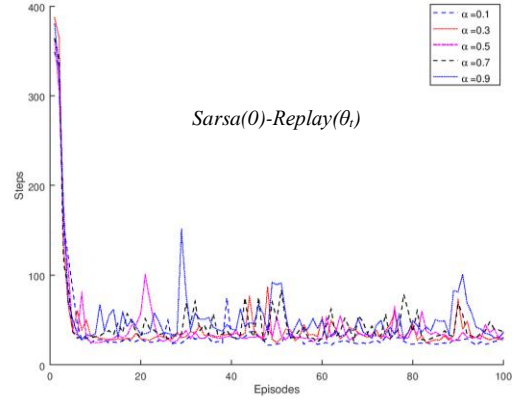


Fig. 5. Comparison of different learning steps for *Sarsa(0)-Replay*( $\theta_t$ ) averaged over 20 runs for binary features.  $\alpha = 0.1$  gave best performance and the algorithm converges very quickly after the 5<sup>th</sup> episode in all cases, lower  $\alpha$  gave more stability for the algorithm. The algorithm is very efficient in terms of complexity although the agent is effectively resamples all of the so far visited states in every time step.

As for the comparison of *Sarsa(0)-Replay*( $\theta_0$ ) and the traditional *Sarsa*( $\lambda$ ). Fig. 6 and Fig 7 show the comparisons results for different learning steps for each algorithm.

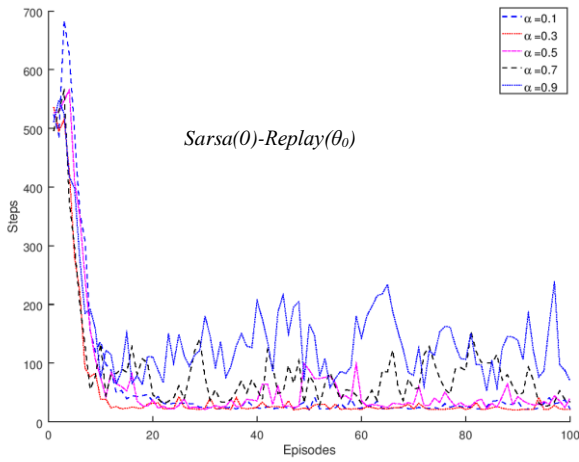


Fig. 6. Compariosn of different learning steps for  $Sarsa(0)-Replay(\theta_{ini})$  averaged over 20 runs for binary features.  $\alpha = 0.1$  is best and the algorithm converges almost as fast as  $Sarsa(0)-Replay(\theta_{ini})$  with even better complexity  $O(n)$ .

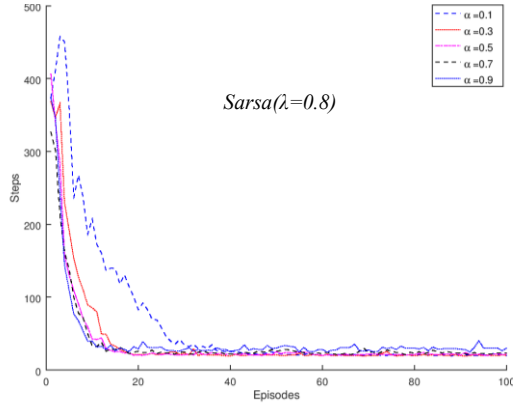


Fig. 7. Compariosn different learning steps for Accumulated  $Sarsa(\lambda=0.8)$  averaged over 20 runs for binary features.  $\alpha = 0.9$  is associated with best perofrmance.

Fig. 8 shows a comparison of TD(0)-Replay along with the other three algorithms. It is again clear that our algorithm exceeds the performance of other algorithms in this domain. It should be noted that  $Sarsa(0)-Replay(\theta_0)$  has the advantage of being very efficient with complexity of  $O(n)$  for each time step calculation, has no hyper paramtrs to be set (such as  $\lambda$ ), it is model free and effetiive for planning and it has an exact online equivalence algorithm that it was build on it; making it also truly online. Fig. 9 shows the execution time for each algorithm. Clearly  $Sarsa(0)-Replay(\theta_t)$  has very good execution time due to its efficiency in converging quickly to an optimal policy. It also can be seen that Dyna Full Planning is the most expensive as expected. Combining execution time along with how fast the algorithm reached the optimal policy, we conclude that  $Sarsa(0)-Replay(\theta_t)$  can be the method of choice for those application that needs to maximise the lived experience, such physical and real time system.

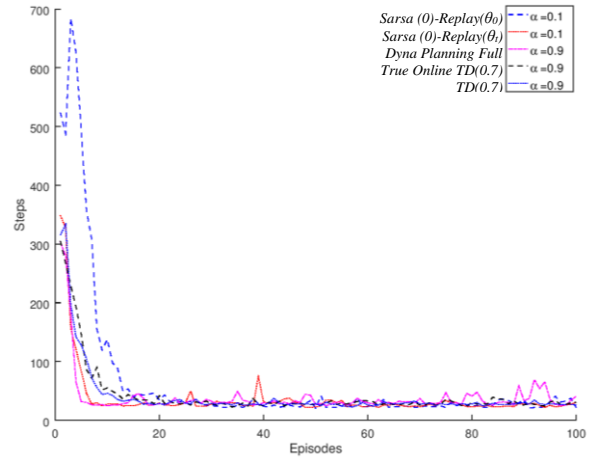


Fig. 8. Final compariosn between  $Sarsa(0)-Replay(\theta_t)$  [complexity is  $O(n \times m)$   $m$  is the numebr of active features], Dyna Full Replanning for Control [complexity is  $O(T^2)$   $T$  is the number of episode step],  $Sarsa(0)-Replay(\theta_0)$  and  $Sarsa(\lambda=0.8)$ .  $Sarsa(0)-Replay(\theta_t)$  has done almost the same job as Dyna Full Re-Planning with much better cost, in fact it was the fastest to converge after 4 episodes. This algorithm is efficient model-free with reaasonable overhead and maximal performance. Next comes  $Sarsa(0)-Replay(\theta_0)$  which has achived very good results compared to the other algorithms including true Online TD with minimal overhead and very good performance. The learning rates have been chosen to maximise the perofrmance of each algorithmt as per the previous set of exeiments,  $\lambda$  is set to 0.8 based on previous studies [8]. It should be noted that better results may have been obtained by setting  $\alpha$  to even a lower values for both TD-Replay algorithms.

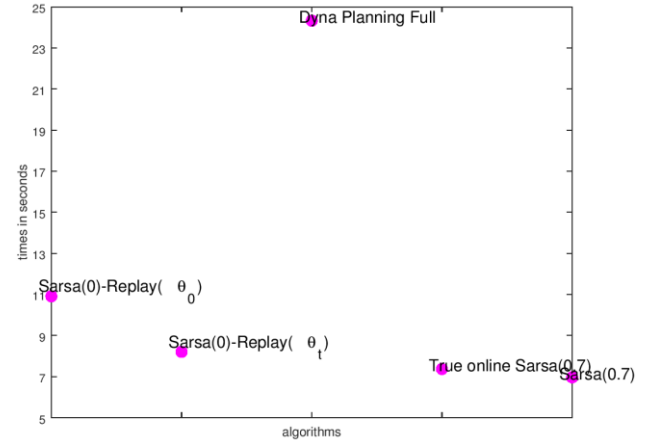


Fig. 9. Compariosn different algorithms execution times.

#### IV. CONCLUSION AND FUTURE WORK

In this paper we have presented three new RL algorithms that allow for an efficient and full replay of all past experience in every step for a reinforcement learning agent life with reasonable and minimal overhead respectively. TD(0)-Replay( $\theta_0$ ) seems good for prediction while  $Sarsa(0)-Replay(\theta_t)$  seems best for control. When we add a delta update for the latter before the replay process the algorithm becomes suitable for planning although it is model-free.



The presented algorithms are suitable for real time and experience-expensive systems (where learning through experience is a hard and expensive process). We showed how to deduce a backward view directly from the forward view for the online case. We have contrasted the presented algorithms with other similar algorithms and showed the differences theoretically in terms of mechanics and practically through experiments. Our experiments confirm the potential for these method to be used in different domains and to overcome other existing replay and planning schemes. In the future we will be looking towards generalising our algorithm to have a general target  $U_i$  instead of the one-step target. This new target will change the form of the matrix  $A$  and special attention would be needed to deduce an efficient form for the full replay of past experience. Another avenue is to develop the algorithms for a non-linear model. Another avenue to develop the algorithms to become a multistep RL algorithms suitable for more general models include averaging [10][11].

## REFERENCES

- [1] Lin, L. J. Self-improving reactive agents based on reinforcement learning, planning and teaching. *Machine Learning*, 8:293–321, 1992.
- [2] Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., et al (2015). Human-level control through deep reinforcement learning. *Nature*, 518:529–533.
- [3] Sutton, R. S., Szepesvari, C., Geramifard, A., and Bowling, M. Dyna-style planning with linear function approximation and prioritized sweeping. In *International Conference on Uncertainty in Artificial Intelligence (UAI)*, pp. 528–536, 2008.
- [4] van Seijen, H. H. and Sutton, R. S. (2015). A Deeper Look at Planning as Learning from Replay. *Proceedings of the 32<sup>nd</sup> International Conference on Machine learning (ICML)*.
- [5] van Seijen, H. H. and Sutton, R. S. (2014). True online TD( $\lambda$ ). In *Proceedings of the 31<sup>st</sup> International Conference on Machine learning (ICML)*.
- [6] Sutton, R. S. and Barto, A. G. (2017). *Reinforcement Learning: An Introduction Complete Draft*. 2nd Edition, Accessed online, MIT Press, Cambridge.
- [7] Bertsekas, D. P. and Tsitsiklis, J. N. (1996). *Neuro-Dynamic Programming*. Athena Scientific.
- [8] Dayan, P. (1992). The convergence of TD( $\lambda$ ) for general  $\lambda$ . *Machine Learning*, 8(3):341–362.
- [9] Sutton, R. S. (1988). Learning to predict by the methods of temporal differences. *Machine Learning*, 3(1):9–4
- [10] van Hasselt, H. and Sutton, R. S. (2015). Learning to predict independent of span. *arXiv:1508.04582*.
- [11] Sutton, R. S. TD models: Modeling the world at a mixture of time scales. In *Proceedings of the 12th International Conference on Machine Learning (ICML)*, pp. 531–539, 1995.